

Predict Used Car Prices

Isaac's R analysis:

Data and Exploratory Analysis of regression of used car prices:

I split the plot window into a 2x2 grid to display boxplots and histograms of the model_year and mileage of vehicles to see the outliers and the distribution of data; model_year is skewed left and mileage is skewed right.

A correlation would make sense here between mileage and model_year because there are more newer cars and more cars with fewer mileage and there are less older cars and less cars with more mileage.

I summarized mileage and model_year to see where the majority of data was.

Model_year summary results:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1974	2013	2017	2016	2020	2024

Mileage summary results:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
100	24500	57500	66043	95798	405000

I stored outliers for mileage and model_year columns so that I could remove them from the data set. I then created a copy of the test data frame and removed the stored outliers from the copy of the data frame.

I summarized mileage and model_year from the new dataset created with outliers removed.

New Model_year summary results:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
2003	2013	2017	2016	2020	2024

New Mileage summary results:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
100	23600	54900	62973	92000	202055

The min for model_year changed but everything else stayed the same, indicating that the outliers were on the lower end of model_year. Everything except for the min changed for the mileage, indicating that the outliers were on the upper end of mileage. I then created new plots of mileage and model_year columns after the outlier removal in the copied data frame. The data became less dramatically skewed but still skewed in the same directions. There were still outliers in the mileage column according to the boxplot, but they weren't outliers according to what math determines to be an outlier.